

PREDICTIVE MODEL AND NEAR INFRARED SPECTROSCOPY IN PREDICTING
THE DIESEL FUEL PROPERTIES

HASAN ALI GAMAL AL-KAF



A project report submitted in partial fulfillment of
the requirement for the award of the Degree of
Master of Electrical Engineering

Faculty of Electrical and Electronic Engineering
Universiti Tun Hussein Onn Malaysia

JANUARY 2018

Dedicate this work to my beloved mother and father and my brothers and my close
friend Ahmed



PTTA UTHM
PERPUSTAKAAN TUNKU TUN AMINAH

ACKNOWLEDGEMENT

First and foremost, I would like to express my gratitude to ALLAH Almighty for His blessing to this great achievement. Without His grace and comparison, none of his would have been possible.

I would like to express my deepest appreciation to my supervisor, Dr Chia Kim Seng for the guidance, enthusiasm for the entire progress of this project.

Special thanks to my family members and friends who have been helpful by offering comment and advice ensuring the success of this research.



PTTA UTHM
PERPUSTAKAAN TUNKU TUN AMINAH

ABSTRACT

Monitoring the diesel fuel properties play an important role in the performance of vehicle engines. Near-infrared (NIR) technology has been investigated as an alternative to monitor the diesel fuel properties. NIR spectroscopy shows an enormous potential for quantitative analysis of complex samples by coupling with artificial neural networks (ANNs). Although a single layer ANN shows promising in the establishing better relationship between a component of interest and NIR spectrum, a different algorithm for updating weight that has been proved to improve the performance of the multilayer could further reveal the potential of single linear layer ANN in NIR spectroscopic analysis. Therefore, this study investigates the performance of a single layer ANN that trained with Levenberg-Marquardt (SLM) and that trained with Scaled Conjugate Gradient (SSCG) and compares the proposed methods with multilayer ANN that trained with same learning algorithms. Results were evaluated and discussed with previous studies that used the same data sets to establish the relationship between the NIR spectral data and diesel fuel properties. Finding depicts that the proposed SLM and SSCG were capable of predicting the diesel fuel properties using NIR spectrum without data reduction, and achieving better accuracy in predicting the diesel fuel properties compared with other recent methods. In addition, using a proposed genetic algorithm for data reduction to improve the predictive model of the proposed method.

ABSTRAK

Pemantauan sifat bahan api diesel memainkan peranan penting dalam prestasi enjin kenderaan. Kajian telah menunjukkan teknologi infra-merah dekat (NIR) boleh digunakan sebagai kaedah alternatif untuk memantau sifat-sifat bahan api diesel. Spektroskopi NIR menunjukkan potensi besar untuk menganalisis kuantitatif sampel yang kompleks dengan gandingan rangkaian saraf tiruan (ANN). Walaupun satu Japisan ANN menjanjikan bubungan yang lebih baik antara komponen berkaitan dan spektrum NIR, algoritma yang berbeza untuk mengemaskini berat telah membuktikan dapat meningkatkan prestasi pelbagai lapisan ANN dapat menunjukkan potensi lapisan linear tunggal ANN dalam analisis spektroskopi NIR. Oleh itu, kajian ini dilakukan untuk menyiasat prestasi lapisan tunggal ANN yang dilatih dengan Levenberg-Marquardt (SLM) dan Gradient Conjugate Scaled (SSCG) dan kemudian dibandingkan dengan kaedah pelbagai lapis ANN yang dilatih dengan algoritma yang sama. Hasil keputusan telah dinilai dan dibincangkan dengan kajian sebelumnya yang menggunakan set data yang sama untuk mewujudkan hubungan antara data spektrum NIR dan bahan api diesel. Penemuan kajian menunjukkan bahawa SLM dan SSCG yang dicadangkan mampu meramalkan sifat bahan api diesel menggunakan spektrum NIR tanpa pengurangan data, dan mencapai ketepatan yang lebih baik dalam meramalkan sifat bahan api diesel berbanding dengan kaedah sebelum ini. Di samping itu, penggunaan algoritma genetik untuk pengurangan data dapat memperbaiki ramalan kaedah menggunakan model yang dicadangkan.

CONTENTS

TITLE	
DECLARATION	11
DICATION	111
ACKNOWLEDGEMENT	iv
ABSTRACT	v
ABSTRAK	vi
CONTENTS	vii
LIST OF TABLES	xi
LIST OF FIGURES	xii
LIST OF APPENDICES	xiii
CHAPTER 1 INTRODUCTION	1
1.1 Project Background	1
1.2 Problem Statemen	3
1.3 Objectives	4
1.4 Scope	5
CHAPTER2 LITERATURE REVIEW	6

2.1	Introduction	6
2.2	Important of diesel fuel properties	6
2.3	Near Infrared Spectroscopy (NIRS)	7
2.4	Multivariable Calibration Analysis	8
2.4.1	Classical Least Squares	9
2.4.2	Inverse Least Square	10
2.4.3	Partial least squares (PLS)	10
2.4.4	Multiple linear regression (MLR)	11
2.4.5	Artificial neural network	12
2.5	Genetic Algorithm	13
2.5.1	Important consideration	14
2.5.2	Genetic algorithm and near infrared spectroscopy	16
2.5.3	Previous Research on Genetic Algorithm with other Regression	16
2.5.4	Neural network coupled with Genetic algorithm	19
2.6	Summary	20
CHAPTER 3	METHODOLOGY	21
3.1	Introduction	21
3.2	Flow of the works	21
3.3	Theory and algorithm	23
3.3.1	Single layer network	23
3.3.2	Multilayer network	24
3.3.3	Levenberg-Marquardt (LM)	25
3.3.4	Scaled Conjugate Gradient (SCG)	25
3.3.4	Single Layer Neural Network	26



3.4	Experimental	27
3.4.1	Efficient of dataset	27
3.4.2	Parameter setting for single layer and multilayer	28
3.5	Genetic algorithm for feature selection	29
3.6	Evaluation the performance	32
3.6.1	Root mean square error	32
3.6.2	Correlation Coefficient	33
CHAPTER4	RESULTS AND DISCUSSION	34
4.1	Introduction	34
4.2	Optimization of iteration of SLM/SCG	34
4.3	Evaluation of the single layer and multilayer trained with lm and SCG	37
4.4	Evaluation of the single layer with previous studies	41
4.5	Genetic algorithm with single layer	43
4.5.1	First Genetic algorithm with single layer trained with SCG	43
4.5.2	First Genetic algorithm with single layer trained with LM	45
4.5.3	Second Genetic algorithm with single layer trained with SCG	45
4.5.4	Second Genetic algorithm with single layer trained with LM	45
CHAPTER5	CONCUSSION AND RECOMMENDATIONS	47
5.1	Conclusion	47
5.2	Recommendation	48



PTT AUTHM
PERPUSTAKAAN TUNKU TUN AMINAH

REFERENCE

49

APPENDIX

56



PTTA UTHM
PERPUSTAKAAN TUNKU TUN AMINAH

LIST OF TABLES

TABLE	TITLE	PAGE
2.1	Previous Research on Genetic Algorithm with other Regression	17
2.2	Selection method, Generation and fitness function of previous research	18
3.1	Example of preparation of chromosomes	29
4.1	Comparison between single layer and multilayer	40
4.2	Comparison between Adaline and proposed methods	41
4.3	Compare the performance of the proposed method and Partial Least Square (PLS) extreme learning machine (ELM) and boosting extreme learning machine (Boosting ELM)	42
4.4	Comparison between proposed method and genetic inverse least squares	44

LIST OF FIGURES

FIGURE NO	TITLE	PAGE
2.1	Objective of Multivariable Regression	8
2.2	Type of Multivariate Calibration Models	9
2.3	Producer of Genetic Algorithm	14
3.1	Overall works of the methodology	22
3.2	General Description of SLM and SSCG models	27
3.3	Flowchart of steps of Genetic algorithm	31
4.1	The performance of the single layer neural network for different properties with different iteration for LM algorithm to find optimized iteration	35
4.2	The performance of the single layer neural network for different properties with different iteration for SCG algorithm to find optimized iteration	36
4.3	The regression plot of the proposed single layer ANN trained by LM for prediction	37
4.4	The regression plot of the proposed single layer ANN trained by SCG for prediction	38
4.5	The performance of the multi-layer ANN that trained by SCG and LM algorithms in predicting the boiling point	39
4.6	The performance of first and second genetic algorithm for different LM and SCG algorithm	46

LIST OF APPENDICES

APPENDIX	TITLE	PAGE
A	Project Gantt chart	56
B	First Genetic Algorithm with Single Layer Trained with SCG	58
C	First Genetic Algorithm with Single Layer Trained with LM	61
D	Second Genetic Algorithm with Single Layer Trained with SCG	64
E	Second Genetic Algorithm with Single Layer Trained with LM	67



PTTA UTHM
PERPUSTAKAAN TUNJUNGAN AMINAH

CHAPTER 1

INTRODUCTION

1.1 Project Background

Near-Infrared (NIR) spectroscopy has been widely proposed as an alternative for the determination of various chemical compounds, e.g. the products of petroleum refining and petrochemicals, food products, and pharmaceutical [1–3]. It has proved its efficiency for laboratory and industrial applications with the advantages of rapid, low cost, and non-invasive. One of the promising usages of NIR spectroscopy is in petroleum industry [4]. This is because petroleum refining and petrochemicals consist of hydrocarbons, which can be quantified by NIR sensing technique. Classical calibration methods i.e. linear and multiple regression is important in helping researchers or users to understand the relationship between the dependent and independent variables. Unfortunately, these classical calibration approaches are unable to directly model complex and high dimensional NIR data. Thus, pre-processing and data reduction strategies are necessary to be performed prior the use of classical calibration approaches so that sophisticated NIR spectral data can be modeled.

NIR spectroscopy shows an enormous potential for quantitative analysis of complex samples by coupling with artificial neural networks (ANNs). Many studies use the NIR spectroscopy coupled with multilayer neural network (MNN) for quantifying the concentrations of urea, creatinine, glucose, oxyhemoglobin [5], the rice wine age [6], the soluble solid content of intact pineapple [7], the fiber contents of textile mixture [8], and the cephalixin [9]. The applicability of the ANN approach has been increased due to the

advances of learning algorithms such as Levenberg-Marquardt, Scaled Conjugate Gradient, Gradient Descent, and One Step Scant.

Nevertheless, multilayer neural network (MNN) is complex in terms of determining its elements i.e. the number of a hidden neuron, the initial weight, and other learning parameters. In most cases, the idea about the right number of hidden neurons might not be clear to identify. Therefore, this leads to many trial-and-error times for finding a right model for an application [10, 11]. Moreover, an excessive number of hidden layers or hidden neuron could negatively affect the generalization of MNN that leads to over-fitting issues [12].

Recently, boosting extreme learning machine (ELM) for single hidden layer feedforward neural networks has been introduced to the local minimum without complex parameters tuning [13]. The boosting ELM has two parameters to be optimized i.e. the activation function and the number of hidden nodes. The inherent characteristics of simple structure, excellent predictive performance, and high learning speed attribute the superiority of boosting ELM [13]. Both weights and biases for hidden nodes are randomly generated without iteratively adjusting in ELM1 [4, 15]. Even though this will increase the learning speed and reduce the number of parameters that need to be optimized, the initialization of the hidden layer biases and input weights that done randomly may make ELM unstable in practices and the complex structure of ELM [15, 16]. Additionally, the complexity of ensemble modeling strategy e.g. the boosting ELM that consists of many sub-models could be a barrier for researchers to understand the fundamental relationship between near infrared wavelengths and the component of interest.

The simplest structure of ANN i.e. single layer ANN that coupled with Levenberg-Marquardt [17] or Gradient Descent technique [18] has shown a potential to achieve faster convergence with lower computational complexity implemented in system identification e.g. real-time adaptive control and online system identification. On the other hand, a single linear layer that trained with Widrow-Hoff delta rule has been successfully implemented to model high dimensional NIR spectral data to predict the boiling point of diesel fuel without any data reduction approach, and achieved better performance compared with principal component regression (PCR) and Partial least square (PLS)[19]. However, the single layer that trained with Widrow-Hoff delta rule depends on two important factors

i.e. the adaption cycle and learning rate. Consequently, a trial and error approach is needed to optimize these two parameters to avoid overfitting problems. Perhaps different algorithms e.g. Levenberg-Marquardt, Scaled Conjugate Gradient for updating weight that has been proved to improve the performance of the multilayer could further reveal the potential of single linear layer ANN in NIR spectroscopic analysis.

In summary, the limitation of the current approaches i.e. the updating weight is a crucial task to improve the accuracy of a single layer ANN. Second, for our best of knowledge different algorithms such as LM and SCG has not yet applied for the single layer ANN in NIR studies. Third, the number of hidden neurons of both MNN and ELM are required to be optimized while single layer ANN e.g. Adaline does not need to optimize this parameter. Fourth, the predictive accuracy is not proportional to the model complexity. This is because a multilayer ANN that gives a good training accuracy may lead to worse generalization than a single layer in some cases. In addition, Spectral data of near infrared consisting of hundreds and even thousands of absorbance values per spectrum. Some of these values are irreverent and not interest and existing of outlying samples and nonlinearity characteristics of NIR spectrum can degrade the performance of a predictive model significantly. Genetic algorithm widely used as features selection to improve the predictive models.

1.2 Problem Statement

Monitoring the diesel fuel properties play an important role in the performance of vehicle engines. Near infrared (NIR) technology has been investigated as an alternative to monitoring the diesel fuel properties. Spectral data of Near infrared consisting of hundreds and even thousands of absorbance values per spectrum. Some of these values are irreverent and not interest and existing of outlying samples and nonlinearity characteristics of NIR spectrum can degrade the performance of a predictive model significantly. Therefore, NIR spectroscopy coupling with artificial neural networks (ANNs) has improved the performance of predicted models. The most widely used is multilayer neural network which shows a good performance but multilayer neural network (MNN) is not a simple model in terms of determining its elements i.e. the number of a hidden neuron, the

initial weight, and other learning parameters. In most cases, the idea about the right number of hidden neurons might not be clear to identify. Therefore, this leads to many trial-and-error times for finding a right model for an application. Moreover, an excessive number of hidden layers or hidden neuron could negatively affect the generalization of MNN that leads to over-fitting issues.

Recently, extreme learning machine (ELM) for single-hidden layer feedforward neural networks (SLFNs) has been introduced to the local minimum without complex parameters tuning [19]. The ELM has two parameters to be optimized i.e. the activation function and the number of hidden nodes. The inherent characteristics of simple structure, excellent predictive performance, and high learning speed attribute the superiority of ELM [19]. Both weights and biases for hidden nodes are randomly generated without iteratively adjusting in ELM. Even though this will increase the learning speed and reduce the number of parameters that need to be optimized, the initialization of the hidden layer biases and input weights that done randomly may make ELM unstable in practices.

Single layer trained with withdraw Hoff (Adaline) is simple structure and has a strong relationship between the near infrared spectrum and the physical properties, but the performance does not give a good performance compared with MNN. In addition, Adaline depends on two important factors i.e. the adaption cycle and learning rate. Consequently, a trial and error approach is needed to optimize these two parameters to avoid overfitting problems. Perhaps different algorithm for updating weight that has been proved to improve the performance of the multilayer could further reveal the potential of single linear layer ANN in NIR spectroscopic analysis.

1.3 Objectives

In this thesis, the overall goal is to predict boiling point of diesel using near infrared spectral data and predictive model. In order to achieve this, the following objectives are listed:

- 1 Develop a single layer ANN and genetic algorithm based ANN in NIR spectroscopic analysis.

- 2 To establish the relationship between the NIR spectrum and the diesel fuel properties using the developed ANN.
- 3 To evaluate the performance of the proposed single layer ANN with multilayer networks that trained with same learning algorithms without data reduction, and the previous studies that used same NIR data.

1.4 Scope

In this thesis, the overall scopes are to provide a clear idea on how genetic algorithm and single layer trained with Levenberg-Marquardt, Scaled Conjugate Gradient is a good technique for optimization the diesel fuel properties by using near infrared spectroscopy.

In order to achieve this, the following scopes are listed below:

- 1 The NIR spectral data of diesel fuel samples that measured at the Southwest Research Institute (SWRI) will use in this study.
- 2 Data processing and modeling will use by MATLAB (version R2015b, win64)
- 3 A single layer trained with Levenberg-Marquardt, Scaled Conjugate Gradient will use as a multivariable technique to predict diesel fuel properties based on near infrared spectrum.
- 4 Multilayer neural network (MNN) is used as regression to compare the accuracy with proposed methods.
- 5 Evaluated and compared with previous studies that used same NIR datasets such as partial least square (PLS), adaptive linear neuron(Adaline), extreme learning machine (ELM), Boosting(ELM) and inverse genetic inverse least squares
- 6 A genetic algorithm is used as feature selection for NIR Spectral data of diesel samples.

CHAPTER 2

LITERATURE REVIEW

2.1 Introduction

In this chapter, the general concept that is related to predicting model using near infrared spectroscopy was reviewed and the significance that motivates in conducting this research was presented in the subsequent chapters of this thesis. First, Section 2.2 depicts the importance of diesel fuel properties. Second, Section 2.3 a brief about near infrared spectroscopy as well as the usage of near infrared spectroscopy in diesel fuel were discussed and presented. Then, Section 2.4 Multivariable techniques were discussed and Section 2.5 focusing on the artificial neural network. Finally, Section 2.6 The Genetic technique include the producer of the genetic algorithm, an important consideration should be taken while using a genetic algorithm and finally the previous research that combines genetic algorithm with other regression were discussed and focusing on an artificial neural network coupled with a genetic algorithm.

2.2 Important of Diesel Fuel Properties

The significance of the different fuel parameters and the ASTM reference method used to determine them are described below. The boiling point, the appearance of components of a high boiling point in fuels can affect the degree of formation of solid combustion products. The freezing temperature is a significant indicator to emphasize

the smooth supply of fuel in an engine and coming to the flowing properties of diesel fuel, the viscosity and density affect these properties in the pipeline. One of the most fundamental physical properties is density which, in conjunction with other properties, can be used to characterize both heavy and light fractions of petroleum and petroleum products. Also, determining the density (or relative density of petroleum and its products) is important for the process to convert the measured volumes to the volumes at the standard temperature of 15 °C. Cetane number is another important property, fuel can be divided into different brands based on the cetane number. In compression ignition engines, the cetane number, also, provides a measure of the ignition characteristics of diesel fuel oil, and is, also, used by petroleum refiners, engine manufacturers, marketers and in commerce as a primary specification measurement related to matching fuels with engines. The sediment in the engine is increased by the aromatics, and toxic substances are produced by aromatic hydrocarbons.

2.3 Near Infrared Spectroscopy (NIRS)

Near-Infrared (NIR) spectroscopy has been widely proposed as an alternative for the determination of various chemical compounds, e.g. the products of petroleum refining and petrochemicals, food products, and pharmaceutical [1–3]. It has proved its efficiency for laboratory and industrial applications with the advantages of rapid, low cost, and non-invasive. One of the promising usages of NIR spectroscopy is in petroleum industry [4]. This is because petroleum refining and petrochemicals consist of hydrocarbons, which can be quantified by NIR sensing technique. Near infrared reflectance spectra (NIRS) has become an effective method for rapid and real-time analysis of fuel properties [20-30]. Organic functional groups that contain hydrogen can form broad absorption bands within near infrared wavelength range. Diesel fuel mainly consists of organic materials such as cetane, which contains a large number of functional groups. Therefore, it is appropriate to apply NIRS to rapid analysis of diesel fuel properties.

Considerable progress has been reported over the last decade [21-29] e.g., Felicio measured the composition of fuel by NIRS [27]; Coope predicted density, viscosity, and boiling point at 50% recovery of diesel fuel using NIRS [28]; Dehui Wu tried online measurement for diesel fuel with NIRS [29].

2.4 Multivariable Calibration Analysis

Multivariate Calibration analysis is very important of NIR data because of the overlapping nature of NIR spectra for condensed phase and the sensitivity of these spectra. These methods are used to achieve robustness and to enhance selectivity by basing measurements on an analysis of the full spectrum. The main objective of the multivariate calibration method is to relate a chemical or physical property of interest to the spectral information encoded across multiple wavelengths by establishes a mathematical model. There are two-step procedures. The first step is to collect a set of standard sample and this sample is pre-determined by an established independent reference assay. Second is to treat the sample spectra and reference as the calibration data for the purpose of establishing a mathematical can correlate the target property to feature in the spectral data set. Finally, the mathematical model can use to predict the model property for subsequent unknown samples. Figure 2.1 shows the producer and the objective of Multivariate Calibration analysis.

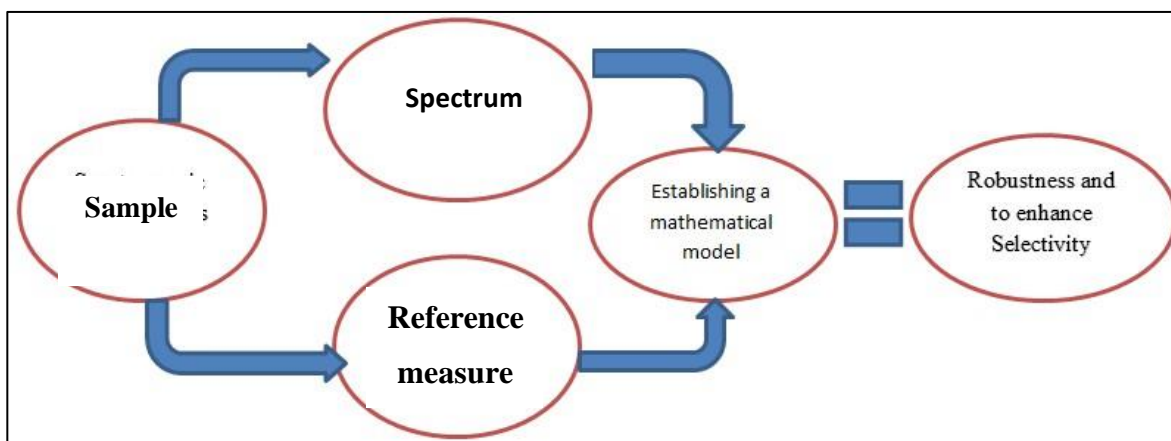


Figure 2.1: Objective of Multivariable Regression

Figure 2.2 shows a variety of multivariate calibration methods are available for the analysis of NIR spectral data. classical least-squares (CLS) and inverse least-squares (ILS) regressions, multiple linear regression (MLR), principal component analysis and regression (PCA and PCR), partial least squares (PLS) and net analyte signal (NAS).

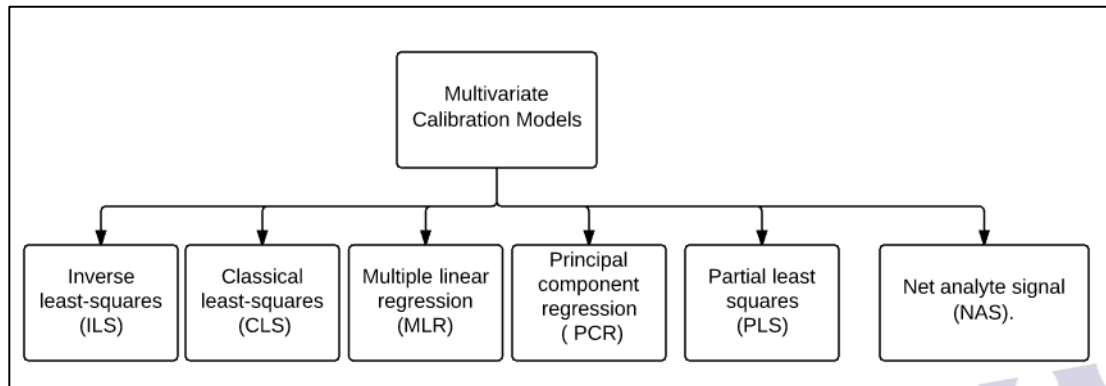


Figure 2.2: Type of Multivariate Calibration Models

2.4.1 Classical Least Squares

Classical Least Squares (CLS) models the spectral variances based on the Beer-Lamber Law. (CLS) and inverse least-squares (ILS) regressions which are often referred as K and P matrix methods, respectively [31-32]. The constituent's concentration and the corresponding are an important parameter as shown in Equation 2.1 [32-33] The constituent spectra \mathbf{S} can be estimated simply with an ordinary least squares pseudo-inversion of \mathbf{C} as shown in Equation 2.2. To predict the concentration of analytes in a new sample from collected spectra is shown in Equation 2.3 and the pseudo-inversion of estimated constituent spectra \mathbf{S} is defined in equation.

$$\mathbf{R} = \mathbf{S}\mathbf{C} + \mathbf{E} \quad (2.1)$$

$$\mathbf{S}^T = \mathbf{R}^T + \mathbf{E}^T \quad (2.2)$$

$$\mathbf{S} = \mathbf{R}^T + \mathbf{E}^T \quad (2.3)$$

Where \mathbf{R} is the vector of sample spectra, \mathbf{S} contains the spectra of each individual constituent in \mathbf{R} and is determined by regression, \mathbf{C} is the corresponding concentration matrix of analyses of interest and \mathbf{E} is the error matrix associated with the calibration.

CLS manipulate the spectra and concentration in an easy way and the constituent spectra estimate through the regression. The disadvantage of this method is the need to have a component spectrum for each element which composites the spectrum. A new approach called advanced CLS (ACLS) could reach the robust modeling within a complex matrix.

2.4.2 Inverse Least Square

The concentration of sample predicted quantitatively from sample spectra [34]. The equation 2.4 shows the basic model for ILS. The vector P is designed to selective to the analyte for modeling and orthogonal to all other components of spectral variation in the spectra matrix R . the regression vector is found in equation 2.5. The advantage of this method is that has a purely mathematical construction.

$$R = RP + E \quad (2.4)$$

$$P = R^+ C \quad (2.5)$$

Where R is the matrix of spectral data, R^+ is the pseudo-inverse of the spectral matrix R .

The vector P is the correlation between the spectral matrix R and the concentration matrix C . E is the matrix of concentration errors that are not included by the model.

2.4.3 Partial Least Squares (PLS)

The Partial least squares were proposed by H.Wold in the 1960s [35]. PLS can estimate the model between spectral variance and concentration in linear regression [36-37]. The basic concept of this method is to decompose the spectral matrix R and the concentration matrix C into loading and score. The second step is to establish the correlation of covariance as shown in equation 2.8. To ensure maximum correlation with the loadings from the concentration, the latent variables from spectral matrix R in PLS are developed as a linear combination of the spectral information.

The advantage of this is to enhance a subspace in the spectral matrix R that is better to model concentration. The relevant equations are:

$$R = R^* + E \quad (2.6)$$

$$C = C^* + F \quad (2.7)$$

Where T and U are score matrices of R and C and, P and R are the respective loading matrices. E and F are the errors associated with the corresponding R and C matrices. the normalized first weight vector which about the connection is between matrices R and C is constructed in order to maximize correlation with the concentration matrix C . the w is computed as follow.

$$w_1 = \frac{R^*{}^T C}{\|R^*{}^T C\|} \quad (2.8)$$

Consequently, the regression vector related to the first factor can be estimated as

$$b_1 = \frac{C^T (R^*{}^T w_1)}{\|(R^*{}^T w_1)^T R^* w_1\|} \quad (2.9)$$

Subsequently, the algorithm repeats for additional factors by calculating w_2 and further regression vectors with the concentration and spectra residuals based on the first factor.

For concentration predictions, the calibration vector matrix is used according to Equation 2.10.

$$\hat{C} = R^* b \quad (2.10)$$

Because the PLS model extracts the latent variables by building covariance between the spectral and concentration matrices, the PLS method efficiently minimizes the impact of interferences. For this reason, PLS methods represent popular tools for multivariate measurements in complex samples.

2.4.4 Multiple Linear Regression (MLR)

This technique used an ordinary least square algorithm to compute the pseudo-inverse of spectra as shown in equation 2.11. The method models a linear relationship between one dependent variable or more [38].

$$R^+ = (R^*{}^T R^*)^{-1} R^*{}^T \quad (2.11)$$

The pseudo-inverse utilized the inversion of $R^*{}^T R^*$, This technique suffers from its instinctive shortcomings and is very sensitive to the conditioning of the spectral matrix.

Estimation of the true inversion of RTR is difficult and inaccurate because of the collinearity existing in the spectral matrix R. Spectral feature are highly overlapped between components and due to this reason, the resulting collinearity in the spectral matrix R diminishes performance of the MLR approach. The number of spectral variables is more than the independent spectra in the matrix R which result in inaccurate estimation of the regression vector. In order to compensate the limitations, a latent variable algorithm such as PCA and PLS have been developed.

2.4.5 Artificial Neural Network

NIR spectroscopy shows an enormous potential for quantitative analysis of complex samples by coupling with artificial neural networks (ANNs). Many studies use the NIR spectroscopy coupled with multilayer neural network (MNN). The applicability of the ANN approach has been increased due to the new advances and recent popularization of learning algorithms such as Levenberg-Marquardt (LM), Scaled Conjugate Gradient (SCG), Gradient Descent (GD), and One Step Scant (OSS). For example, ANN has been investigated and implemented in predicting the pre-oxidation efficiency of refractory gold concentrate [39] the blood glucose [40]. The emissions from swine buildings [41], the software defect [42], the risks of capital flow [43], and the sleep disorders [44]. Nevertheless, Multilayer neural network (MNN) is not a simple model in terms of determining its elements i.e. the number of a hidden neuron, the initial weight, and other learning parameters. In most cases, the idea about the right number of hidden neurons might not be clear to identify. Therefore, this leads to many trial-and-error times for finding a right model for an application [10 11]. Moreover, an excessive number of hidden layers or hidden neuron could negatively affect the generalization of MNN that leads to over-fitting issues [12].

Recently, extreme learning machine (ELM) for single-hidden layer feedforward neural networks (SLFNs) has been introduced to the local minimum without complex parameters tuning [13]. The ELM has two parameters to be optimized i.e. the activation function and the number of hidden nodes. The inherent characteristics of simple structure, excellent predictive performance, and high learning speed attribute the superiority of ELM

REFERENCES

- [1] D. Wu, X. Chen, P. Shi, S. Wang, F. Feng, and Y. He, "Determination of α -linolenic acid and linoleic acid in edible oils using near-infrared spectroscopy improved by wavelet transform and uninformative variable elimination", *Analytica Chimica Acta*, vol. 634, no. 2, pp. 166-171, 2009.
- [2] D. Wu, Y. He, P. Nie, F. Cao and Y. Bao, "Hybrid variable selection invisible and near-infrared spectral analysis for non-invasive quality determination of grape juice", *Analytica Chimica Acta*, vol. 659, no. 1-2, pp. 229-237, 2010.
- [3] D. Wu, Y. He, J. Shi and S. Feng, "Exploring Near and Midinfrared Spectroscopy to Predict Trace Iron and Zinc Contents in Powdered Milk", *Journal of Agricultural and Food Chemistry*, vol. 57, no. 5, pp. 1697-1704, 2009.
- [4] R. Balabin, R. Safieva and E. Lomakina, "Gasoline classification using near infrared (NIR) spectroscopy data: Comparison of multivariate techniques", *Analytica Chimica Acta*, vol. 671, no. 1-2, pp. 27-35, 2010.
- [5] D. Kalamatianos, P. Liatsis, and P. E. Wellstead, "Near-infrared spectroscopic measurements of blood analytes using multi-layer perceptron neural networks.," *Conf. Proc. IEEE Eng. Med. Biol. Soc.*, vol. 1, pp. 3541–3544, 2006.
- [6] F. Liu, F. Cao, L. Wang, and Y. He, "Discrimination of Rice Wine Age Using Visible and Near Infrared Spectroscopy Combined with BP Neural Network, In *Image and Signal Processing*, pp. 267–271, 2008.
- [7] K. S. Chia, H. A. Rahim, and R. A. Rahim, "Artificial Neural Network Coupled with Robust Principal Components in Near Infrared Spectroscopic Analysis," In *Signal Processing and its Applications (CSPA)*, pp. 19–22, 2012.
- [8] L. Liu, L. Yan, Y. Xie, and G. Xia, "Content Measurement of Textile Mixture by Near Infrared Spectroscopy Based on BP Neural Network," In *Image and Signal*

- Processing (CISP), 2010 3rd International Congress on, vol. 7, pp. 3354–3358, 2010.
- [9] Q. Fei, M. Li, B. Wang, Y. Huan, G. Feng, and Y. Ren, "Analysis of cefalexin With NIR spectrometry coupled to artificial neural networks with modified Genetic algorithm for wavelength selection", *Chemometrics and Intelligent Laboratory Systems*, vol. 97, no. 2, pp. 127-131, 2009.
- [10] S. R. B, "Role of Hidden Neurons in an Elman Recurrent Neural Network in Classification of Cavitation Signals," *Int J Comput Appl*, vol. 37, no. 7, pp. 9–13, 2012.
- [11] Z. Yu-xin and W. Hao-yu, "Intrusive Detection Systems Design based on BP Neural Network," In *Distributed Computing and Applications to Business Engineering and Science (DCABES)*, 2010 Ninth International Symposium on, pp. 462–465, 2010.
- [12] Bouzida Y, Cuppens F."Neural networks vs. decision trees for intrusion Detection". *Commun 2006 ICC '06 IEEE Int Conf.* 2006;2394–400.
- [13] Bian X, Zhang C, Tan X, Dymek M, Guo Y, Lin L, et al. "A boosting extreme learning machine for near-infrared spectral quantitative analysis of diesel fuel and edible blend oil samples". *Anal Methods*. 2017;9(20):2983–9.
- [14] Bian X-H, Li S-J, Fan M-R, Guo Y-G, Chang N, Wang J-J."Spectral quantitative analysis of complex samples based on the extreme learning machine". *Anal Methods*. 2016;8(23):4674–9.
- [15] Chen W-R, Bin J, Lu H-M, Zhang Z-M, Liang Y-Z. "Calibration transfer via an extreme learning machine auto-encoder". *Analyst*. 2016;141(6):1973–80.
- [16] Lu H, An C, Zheng E, Lu Y. "Dissimilarity based ensemble of extreme learning machine for gene expression data classification". *Neurocomputing*. 2014;128:22–30.
- [17] W. Zhang and L. Rock, "An Extended ADALINE Neural Network Trained by Levenberg- Marquardt Method for System Identification of Linear Systems," In *Control and Decision Conference (CCDC)*, 2013 25th Chinese pp. 2453– 2458, 2013.
- [18] S. Bhama and H. Singh, "Single layer neural networks for linear system

- identification using\ngradient descent technique,” IEEE Trans. Neural Networks, vol. 4, no. 5, pp. 884–888, 1993.
- [19] K. S. Chia, “Predicting the Boiling Point of Diesel Fuel using Adaptive Linear Neuron and Near Infrared Spectrum,” In Control Conference (ASCC), 10th Asian pp. 0–2, 2015.
- [20] G. Knothe, “Determining the blend level of mixtures of bio-diesel with conventional diesel fuel by fiber-optic near-infrared spectroscopy a nuclear magnetic resonance spectroscopy”, J. Am. Oil. Chem. Soc. 78 (2001) 1025–1028.
- [21] B.L.L. de Fatima, F.V.C. De Yasconcelos, C.F. Pereira, et al., “Prediction of properties of diesel/biodiesel blends by infrared spectroscopy and multivariate calibration”, Fuel 89 (2010) 405–409.
- [22] I.K. de Oliveira, R.W.F. de Carvalho, R.J. Poppi, “Application of near infrared spectroscopy and multivariate control charts for monitoring biodiesel blends”, Anal. Chim. Acta 642 (1) (2009) 217–221.
- [23] W.J. Gao, H.F. Yuan, X.Y. Li, “Online analysis of heavy alkyl benzene using near infrared spectroscopy, Petrochem”. Technol. 39 (2010) 1388–1389.
- [24] C.L.A. Ulio, B.H. Claudete, J.P. Ronei, “Determination of diesel quality Parameter using Support vector regression and near infrared spectroscopy for an in-line blending optimizer System”, Fuel 97 (2012) 712–715.
- [25] W.B. Zhang, W.Q. Yuan, X.M. Zhang, “Predicting the dynamic and kinematic viscosities of biodiesel blends using mid and near infrared spectroscopy”, Appl. Energy 98 (2012) 123–125.
- [26] L.M. Fang, M. Lin, “Near infrared spectroscopy analysis of diesel fuel by independent component analysis”, Acta Petrolei Sinica (Petroleum Process. Sect.). 24 (2008) 729–731.
- [27] C.C. Felieio, L.P. Bras, J.A. LoPes, “Comparison of PLS algorithms in gasoline And gas oil parameter monitoring with MIR and NIR, Chemome”. Intell. Lab. Syst. 78 (2005) 74–78.
- [28] J.B. Cooper, C.M. Larkin, J. Sehmitigal, et al., “Rapid analysis of jet fuel using a

- hand-held near-infrared (NIR) analyzer”, *Appl. Spectroscopy*. 65 (2011) 187–192.
- [29] D.H. Hu, “Online analysis for diesel based on near infrared spectroscopy, *Spectroscopy*”. *Spectral Anal.* 28 (2008) 1530–1534.
- [30] Y.S. Wang, M. Yang, G. Wei, et al., “Improved PLS regression based on SVM classification for rapid analysis of coal properties by near-infrared reflectance spectroscopy”, *Sens. Actuators B: Chem.* 193 (2014) 723–729.
- [31] Antoon, M.K.; Koenig, J.H; Koenig, J.L., “Least-squares Curve-fitting of Fourier Transform Infrared Spectra with Applications to Polymer Systems, *Applied Spectroscopy*”, 1977, 31, 518-524.
- [32] Dousseau, F.; Pezolet, M., “Determination of the Secondary Structure Content of Proteins in Aqueous Solution from Their Amide I and Amide II Infrared Bands. Comparison between Classical and Partial Least-squares Methods”, *Biochemistry*, 1990, 29, 8771-8779.
- [33] Haaland, D.M.; Easterling, R.G.; Vopicka, D.A., “Multivariate Least-squares Methods Applied to the Quantitative Spectral Analysis of Multicomponent Samples”, *Applied Spectroscopy*, 1985, 39, 73-84.
- [34] Krutchkoff, R.G., “Classical and Inverse Regression Methods of Calibration”, *Technometrics*, 1967, 9, 425-439.
- [35] Wold, H., “Estimation of Principal Components and Related Models by Iterative Least Squares”. In *Multivariate Analysis*; Krishnaiah P.R., Ed.; Academic Press: New York, 1973, P. 383-407.
- [36] Geladi, P.; Kowalski, B.R., “Partial Least-squares Regression: a Tutorial”, *Analytica Chimica Acta*, 1986, 185, 1-17.
- [37] Haaland, D.M., Thomas, E.V., “Partial Least-squares Methods for Spectral Analyses. I. Relation to Other Quantitative Calibration Methods and the Extraction of Qualitative Information”, *Analytical Chemistry*, 1988, 60, 1193-1202.
- [38] Andrews, D.F., “A Robust Method for Multiple Linear Regression”, *Technometrics*, 1974, 16, 523-531.
- [39] Li QC, Li DX, Chen QY. “Prediction of pre-oxidation efficiency of refractory

- gold concentrate by ozone in ferric sulfate solution using artificial neural networks”. Transactions of Nonferrous Metals Society of China. 2011 Feb 1;21(2):413-22.
- [40] Shanthi S, Balamurugan P, Kumar D. “Performance comparison of featured neural network with gradient descent and levenberg-marquart algorithm trained neural networks for prediction of blood glucose values with continuous glucose monitoring sensor data”. InEmerging Trends in Science, Engineering and Technology (INCOSET), 2012 International Conference on 2012 Dec 13 (pp. 385-391). IEEE.
- [41] Sun G, Hoff SJ, Zelle BC, Smith MA. “Development and comparison of backpropagation and generalized regression neural network models to predict diurnal and seasonal gas and PM10 concentrations and emissions from swine buildings”. In2008 Providence, Rhode Island, June 29–July 2, 2008 2008 (p. 1). American Society of Agricultural and Biological Engineers.
- [42] Arora I, Saha A. “Comparison of back propagation training algorithms for software defect prediction”. InContemporary Computing and Informatics (IC3I), 2016 2nd International Conference on 2016 Dec 14 (pp. 51-58). IEEE.
- [43] Wang XP, Huang YS. “Predicting risks of capital flow using artificial neural network and levenberg marquardt algorithm”. InMachine Learning and Cybernetics, 2008 International Conference on 2008 Jul 12 (Vol. 3, pp. 1353-1357). IEEE.
- [44] Garg VK, Bansal RK. “Comparison of neural network back propagation algorithms for early detection of sleep disorders”. InComputer Engineering and Applications (ICACEA), 2015 International Conference on Advances in 2015 Mar 19 (pp. 71-75). IEEE.
- [45] R. Leardi and A. Lubiáñez Gonzalez, “Genetic algorithms applied to feature selection in PLS regression: how and when to use them”, Chemometr. Intell. Lab. Syst. 41, 195–207 (1998).
- [46] H.C. Goicoechea and A.C. Olivieri, “A new family of genetic algorithms for wavelength interval selection in multivariate analytical spectroscopy”, Chemometrics 17, 338-345 (2003).

- [47] H. C. Goicoechea and A. C. Olivieri, "Wavelength selection for multivariate calibration using a genetic algorithm: A novel initialization strategy," *J. Chem. Inf. Comput. Sci.*, vol. 42, no. 5, pp. 1146–1153, 2002.
- [48] R. M. Jarvis and R. Goodacre, "Genetic algorithm optimization for pre-processing and variable selection of spectroscopic data," *Bioinformatics*, vol. 21, no. 7, pp. 860–868, 2005.
- [49] L. Weiling, L. Libing, C. Yingshu, and Y. Qilian, "Application of GA-PLS for Feature Selection and Calibration in Biological Sample," pp. 432–436.
- [50] A. S. Barros and D. N. Rutledge, "Genetic algorithm applied to the selection of principal components," *Chemom. Intell. Lab. Syst.*, vol. 40, no. 1, pp. 65–81, 1998.
- [51] R. Leardi, M. B. Seasholtz, and R. J. Pell, "Variable selection for multivariate calibration using a genetic algorithm: Prediction of additive concentrations in polymer films from Fourier transform-infrared spectral data," *Anal. Chim. Acta*, vol. 461, no. 2, pp. 189–200, 2002.
- [52] L. Sratthaphut and N. Ruangwises, "Genetic algorithms-based approach for wavelength selection in spectrophotometric determination of vitamin B12 in pharmaceutical tablets by partial least-squares," *Procedia Eng.*, vol. 32, pp. 225–231, 2012.
- [53] Q. Fei, M. Li, B. Wang, Y. Huan, G. Feng, and Y. Ren, "Analysis of cefalexin with NIR spectrometry coupled to artificial neural networks with modified genetic algorithm for wavelength selection," *Chemom. Intell. Lab. Syst.*, vol. 97, no. 2, pp. 127–131, 2009.
- [54] D. I. Broadhurst, R. Goodacre, A. Jones, J. J. Rowland, and D. B. Kell, "Genetic algorithms as a method for variable selection in multiple linear regression and partial least squares regression, with applications to pyrolysis mass spectrometry", *Anal. Chim. Acta*, vol. 348, no. 1–3, pp. 71–86, 1997.
- [55] N. Goldshleger, A. Chudnovsky, and E. Ben-Dor, "Using reflectance spectroscopy and artificial neural network to assess water infiltration rate into the soil profile", *Appl. Environ. Soil Sci.*, vol. 2012, pp. 7–9, 2012.
- [56] GC MP, Shettigar AK, Krishna P, Parappagoudar MB. "Back Propagation Genetic

and Recurrent Neural Network Applications in Modelling and Analysis of Squeeze Casting Process”. *Applied Soft Computing*. 2017 Jun 13.

- [57] Silalahi DD, Reaño CE, Lansigan FP, Panopio RG, Bantayan NC. “Using genetic algorithm neural network on near infrared spectral data for ripeness grading of oil palm (*Elaeis guineensis* Jacq) fresh fruit”. *Information Processing in Agriculture*. 2016 Dec 31;3(4):252-61.
- [58] Fülöp A, Hancsók J. “Comparison of calibration models based on near infrared spectroscopy data for the determination of plant oil properties”. *Hungarian Journal of Industry and Chemistry*. 2009 Sep 1;37(2).



PTTA UTHM
PERPUSTAKAAN TUNKU TUN AMINAH